



# Encounter Complexes For Clustering Network Flow

Leigh Metcalf, lbmetcalf@cert.org  
Flocon 2015  
Date



Report Documentation Page				Form Approved OMB No. 0704-0188	
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE <b>JAN 2015</b>		2. REPORT TYPE		3. DATES COVERED <b>00-00-2015 to 00-00-2015</b>	
4. TITLE AND SUBTITLE <b>Encounter Complexes For Clustering Network Flow</b>				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>Carnegie Mellon University, Software Engineering Institute, Pittsburgh, PA, 15213</b>				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release; distribution unlimited</b>					
13. SUPPLEMENTARY NOTES <b>FloCon 2015, Portland, OR, January 12-15, 2015.</b>					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT <b>Same as Report (SAR)</b>	18. NUMBER OF PAGES <b>37</b>	19a. NAME OF RESPONSIBLE PERSON
a REPORT <b>unclassified</b>	b ABSTRACT <b>unclassified</b>	c THIS PAGE <b>unclassified</b>			

---

Copyright 2014 Carnegie Mellon University

This material is based upon work funded and supported by Department of Homeland Security under Contract No. FA8721-05-C-0003 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center sponsored by the United States Department of Defense.

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

This material has been approved for public release and unlimited distribution except as restricted below.

This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other use. Requests for permission should be directed to the Software Engineering Institute at [permission@sei.cmu.edu](mailto:permission@sei.cmu.edu).

FloCon® is registered in the U.S. Patent and Trademark Office by Carnegie Mellon University.

DM-0001965

# Network Flow

---

Clustering Network Flow for fun and profit!

Previously done for finding Trojans, Botnets, Spoofed flows...

But those methods use 'known behavior' to find repeats of that behavior.

# Network Flow

---

The Encounter Complex uses no prior knowledge in its creation.

It is based on encounter traces, which occur when two nodes meet. We record the time and analyze the data.

Encounter traces can include:

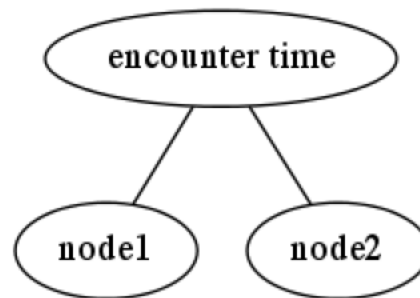
- Two animals meet at a watering hole.

- Two users use the same wireless node.

# Encounter Trace

---

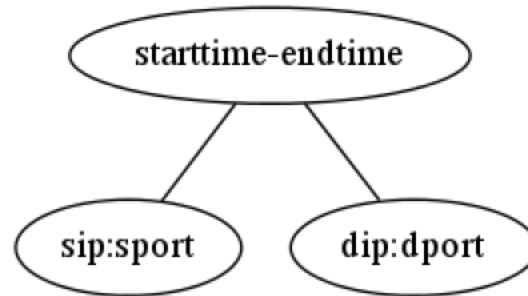
Encounter traces are defined as:



# Encounter Trace for Network Flow

---

Defined as:



I maintain the time period rather than just a single moment in time.

# Encounter Complex

---

Two traces have an edge between them if:

1. They share a node in common
2. The end of one occurs within  $\Delta$  seconds of the start of the next



# Encounter Complex

---

Let's assume  $\Delta=8$

sIP:sPort	dIP:dPort	stime	etime
192.0.2.5:80	192.0.2.200:5265	1412870783	1412870880
192.0.2.199:5353	192.0.2.5:80	1412870885	1412871150

These two flows are connected since the first ends within 5 seconds of the second beginning.

# Encounter Complex

---

Still assuming  $\Delta=8$

sIP:sPort	dIP:dPort	sTime	eTime
192.0.2.5:80	192.0.2.200:5265	1412870783	1412870880
192.0.2.199:5353	192.0.2.5:80	1412870885	1412871150
192.0.2.150:5353	192.0.2.3:25	1412870887	1412871175
192.0.2.5:80	192.0.2.205:5031	1412871160	1412871200

The third row does not share a node in common with the first two.

The second fails  $\Delta=8$  test, but would be part of the complex if  $\Delta \geq 10$

# Encounter Complex

---

We denote the Encounter Complex by  $G_{\Delta}$

Proposition:

$$\text{If } \Delta \leq \Gamma \text{ then } G_{\Delta} \subseteq G_{\Gamma}$$

This is clear because if two nodes are within  $\Delta$  seconds of each other they are certainly within  $\Gamma$  seconds of each other.

# Encounter Complex – Example

---

I used the LBNL data set.

- 11Gb of anonymized data
- Collected from October 2004 through January 2005
- Contains approximately 2.2 million flows
- Covers a wide variety of enterprise traffic

# Encounter Complexes – Example

---

The time I chose had data from two sensors and contained:

- 47,834 network flows
- 1,423 IP addresses
- Average length of flow was 41.34 seconds
- Covered a little over an hour of traffic

# Encounter Complex – Example

---

I created complexes for 7 values of  $\Delta$ :

$\Delta$	Number of Graphs	Edges	Vertices
1	6115	182,485	37,184
50	1498	3,681,789	40,623
100	891	6,769,521	40,763
200	695	12,551,635	40,807
300	597	18,325,825	40,822
400	537	23,605,755	40,831
Infinity	363	106,281,681	40,858

# Encounter Complexes – Example

---

When  $\Delta = \text{infinity}$ , there is quite a lot of work to be done creating the graph. It's essentially  $n^2$  where  $n$  is the number of flows.

It does contain all of the other graphs though...

# Encounter Complexes – Example

---

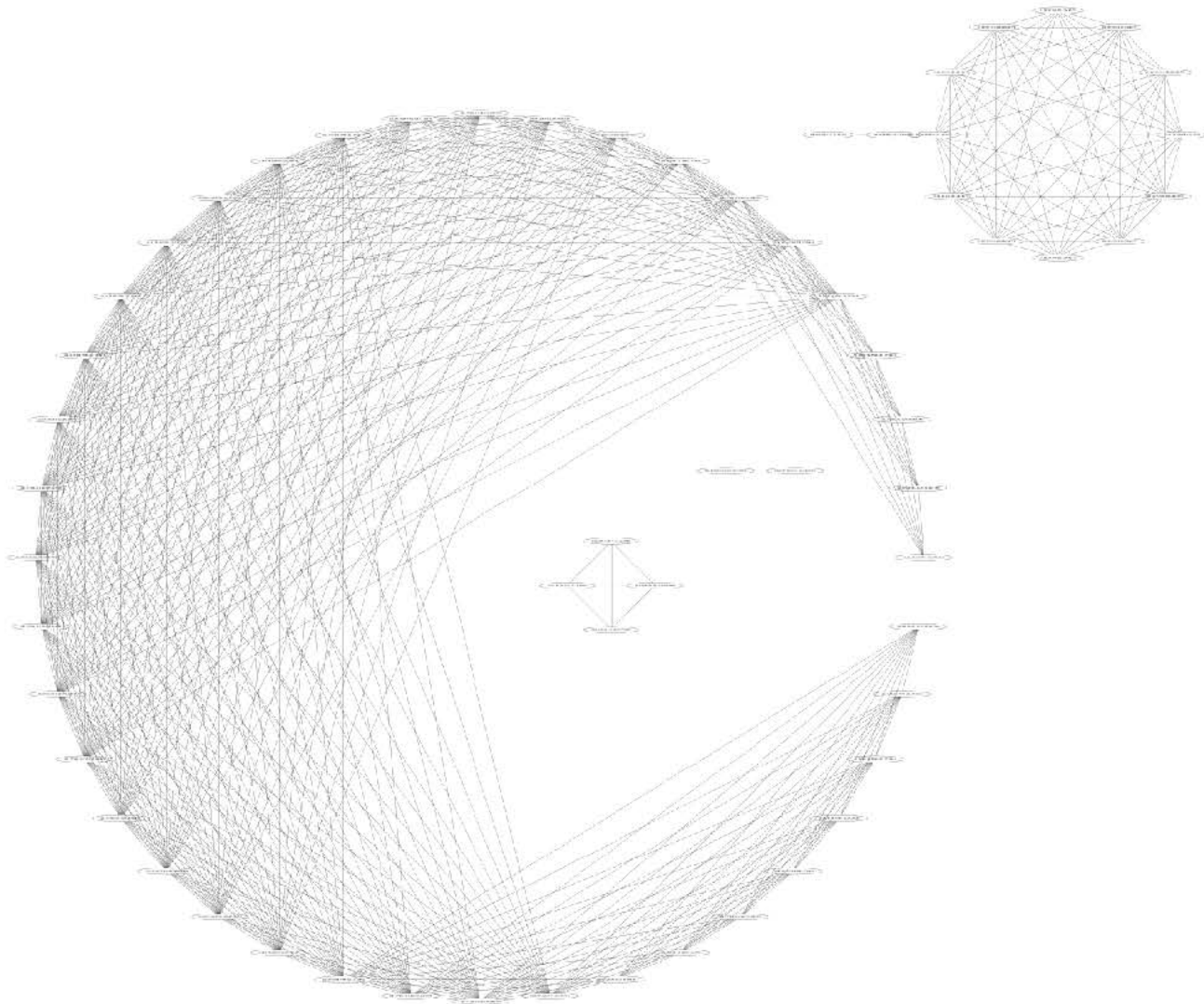
Analyzing it by visualization isn't very useful.

When  $\Delta=1$  there are 6115 graphs to analyze...

Four of which are on the next slide.



# Encounter Complexes – Example



# Encounter Complexes -- Example

---

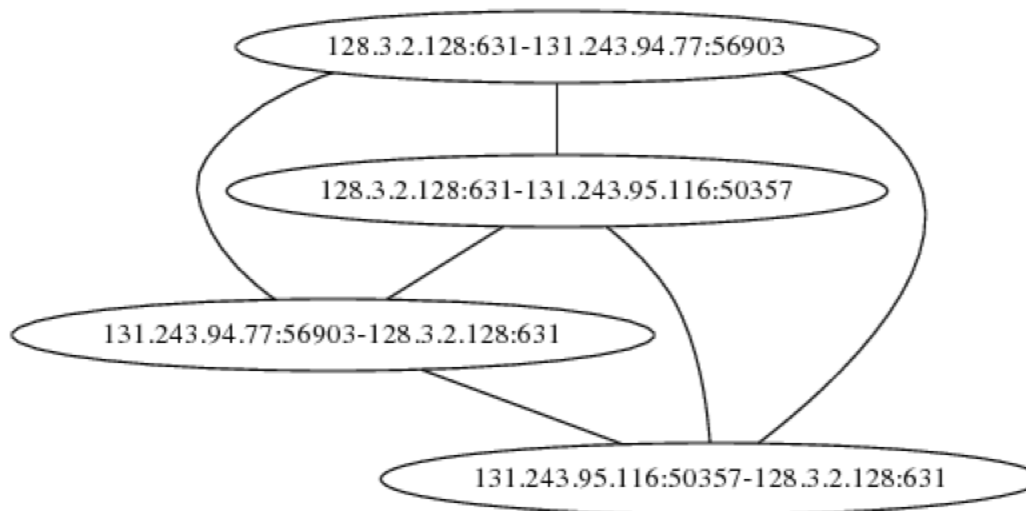
We can analyze them by looking at two things:

- The vertex with the highest degree
- The node that is most prevalent through the graph

# Encounter Complexes -- Example

---

This graph has 128.3.2.128:631 as the most common node.... Could be printing!



# Clustering the Clusters

---

We've created graphs from the network flow...

...now we want to cluster those graphs

- Similar traffic!
- Fewer things to look at!
- Everyone wins!

# Clustering the Clusters

---

There are three steps to creating the clusters:

- 1) Group together those graphs with similar port  
Look at the vertex with the highest degree and consider the ports there.
- 2) Now refine those clusters by putting together graphs with a similar amount of vertices  
Where 'similar amount' means within 10%

# Digression into Graph Theory

---

Graph isomorphisms are an NP complete problem

(This means it is impossible in a reasonable amount of time)

Graph similarity has as many methods as mathematicians working on the problem

...so of course, I came up with my own method.

# Digression into Graph Theory

---

The degrees of vertices within the graph in an encounter complex are a measure of similarity within that graph

The higher the degree, the more similar the vertex is to other vertices in the graph

# Digression into Graph Theory

---

Method:

Given two graphs  $G_1$  and  $G_2$  create a sorted degree vector for each graph. (That is, put all of the degrees of each graph in a vector then sort it.)

If one vector is shorter than the other, pad that one with zeroes until they match in size.



# Digression into Graph Theory

---

We can have two graphs with the same number of edges, vertices and cycles that have different degree vectors.

Example: A graph with 7 edges, 6 vertices had 2 cycles.

$[3, 3, 2, 2, 2, 2]$

$[5, 3, 2, 2, 1, 1]$

Both valid degree vectors for this graph.

# Digression into Graph Theory

---

Once you have the two vectors, use the Pearson coefficient as a distance measure.

Pearson measures the linear dependence between the two vectors.

# Digression into Graph Theory

---

It is also possible to have two graphs that are distinctly different but the Pearson coefficient of the degree vectors is 1.

# Digression into Graph Theory

---

Example:

A graph with 42 edges, 10 vertices and 33 cycles:

[9, 9, 9, 9, 9, 9, 8, 8, 7, 7]

A graph with 29 edges, 10 vertices and 20 cycles:

[7, 7, 7, 7, 7, 7, 5, 5, 3, 3].

Pearson coefficient is 1 in this case.

These two graphs are modelling similar behavior

# Clustering the Clusters

---

Last step of refinement:

- 3) Two graphs are in the same cluster if their Pearson coefficient is greater than 0.9

# Encounter Complexes – Example

---

$\Delta$	Clusters	Number of Clustered Graphs
1	99	756
50	29	193
100	6	32
200	4	11
300	6	14
400	3	6
infinity	0	0

# Encounter Complexes -- Example

---

For  $\Delta=1$  I found a cluster with 47 graphs.

All of these graphs had sIP:50122 in common.

50122 can be used for:

SAP, Symantec and SSH forwarding

Without more information, I don't know much... other than they have common activity across the graphs

# Encounter Complexes – Example

---

Another cluster had 50 graphs.

- Port 80 is common across the cluster
- But no common node

We found similar web traffic patterns



# Comparing Encounter Complexes

---

Clustering the clusters works well when looking at a single complex...

...What if I compare two complexes?

# Comparing Encounter Complexes

---

I chose a second time period from the LBNL data.

Contained:

- 127,223 flows
- 4,490 IP addresses
- A little over an hour of data

I created a complex where  $\Delta=1$

- Contained 14,676 components

# Comparing Encounter Complexes

---

I then compared the two complexes using the criteria listed before:

1. Similar port
2. Similar size
3. Pearson measurement of degree vectors  $> 0.9$

# Comparing Encounter Complexes

---

I found 63 clusters when I compared the two graphs containing a total of 2,087 subgraphs.

I examined one cluster that contained 8 subgraphs

- 2 from one encounter complex

- 6 from the other encounter complex

The common port was 427 but the destination IP address varied

# Future Work

---

- Bytes!
  - Weight the encounter complexes with the bytes transferred in the process
- Protocol!
  - Label the encounter complexes using the protocols in the flow
- Persistent Homology
  - Apply this to the infinity graphs to compare encounter complexes



**Questions/comments?**

